

# A refined characterization of large-scale genomic differences in the first complete human genome

Xiangyu Yang<sup>1,†</sup>, Xuankai Wang<sup>1,†</sup>, Yawen Zou<sup>1</sup>, Shilong Zhang<sup>1</sup>, Manying Xia<sup>1</sup>, Mitchell R. Vollger<sup>2</sup>, Nae-Chyun Chen<sup>3</sup>, Dylan J. Taylor<sup>4</sup>, William T. Harvey<sup>2</sup>, Glennis A. Logsdon<sup>2</sup>, Dan Meng<sup>1</sup>, Junfeng Shi<sup>5,6</sup>, Rajiv C. McCoy<sup>4</sup>, Michael C. Schatz<sup>3,4</sup>, Weidong Li<sup>1</sup>, Evan E. Eichler<sup>2,7</sup>, Qing Lu<sup>1</sup>, Yafei Mao<sup>1,6,\*</sup>

<sup>1</sup> *Bio-X Institutes, Key Laboratory for the Genetics of Developmental and Neuropsychiatric Disorders, Ministry of Education, Shanghai Jiao Tong University, Shanghai, China.*

<sup>2</sup> *Department of Genome Sciences, University of Washington School of Medicine, Seattle, WA, USA.*

<sup>3</sup> *Department of Computer Science, Johns Hopkins University, Baltimore, MD, USA.*

<sup>4</sup> *Department of Biology, Johns Hopkins University, Baltimore, MD, USA.*

<sup>5</sup> *Shanghai Engineering Research Center of Advanced Dental Technology and Materials, Shanghai, China.*

<sup>6</sup> *Shanghai Key Laboratory of Stomatology, Shanghai Ninth People's Hospital, College of Stomatology, Shanghai Jiao Tong University School of Medicine, Shanghai, China.*

<sup>7</sup> *Howard Hughes Medical Institute, University of Washington, Seattle, WA, USA.*

<sup>†</sup> : These authors contributed equally: Xiangyu Yang, Xuankai Wang

\*Corresponding author: [yafmao@sjtu.edu.cn](mailto:yafmao@sjtu.edu.cn)

## Abstract

The first telomere-to-telomere (T2T) human genome assembly (T2T-CHM13) release was a milestone in human genomics. The T2T-CHM13 genome assembly extends our understanding of telomeres, centromeres, segmental duplication, and other complex regions. The current human genome reference (GRCh38) has been widely used in various human genomic studies. However, the large-scale genomic differences between these two important genome assemblies are not characterized in detail yet. Here, we identify 590 discrepant regions (~226 Mbp) in total. In addition to the previously reported ‘non-syntenic’ regions, we identify 67 additional large-scale discrepant regions and precisely categorize them into four structural types with a newly developed website tool (SynPlotter). The discrepant regions (~20.4 Mbp) excluding telomeric and centromeric regions are highly structurally polymorphic in humans, where copy number variation are likely associated with various human disease and disease susceptibility, such as immune and neurodevelopmental disorders. The analyses of a newly identified discrepant region—the *KLRC* gene cluster—shows that the depletion of *KLRC2* by a single deletion event is associated with natural killer cell differentiation in ~20% of humans. Meanwhile, the rapid amino acid replacements within *KLRC3* is consistent with the action of natural selection during primate evolution. Our study furthers our understanding of the large-scale structural variation differences between these two crucial human reference genomes and future interpretation of studies of human genetic variation.

## Introduction

The first draft human genome published two decades ago has contributed enormously to human genomics, medical genomics, evolutionary genomics, and other fields<sup>1,2</sup>. Given efforts to refine and construct the sequence from Genome Reference Consortium, the current human reference genome assembly (GRCh38) has been widely used for understanding human diversity, disease-related variants, and human/primate evolution<sup>3</sup>. The GRCh38 genome assembly has been annotated with abundant resources including gene annotation, gene expression, gene regulation, and others<sup>3</sup>. Despite the high quality of the GRCh38 reference, it still has hundreds of gaps and errors in GRCh38<sup>4</sup>. These gaps and errors represented long-standing obstacles to fully understanding human genomics, especially in repetitive regions<sup>4-10</sup>. With advances in long-read sequencing and computational algorithms, the Telomere to Telomere (T2T) Consortium has finally achieved the goal of building a gapless and accurate assembly of a human genome<sup>4-9</sup>.

The release of the complete genome (T2T-CHM13) provides the first complete sequence view of centromeres, telomeres, tandem repeat arrays, segmental duplications (segdups), and the p-arms of acrocentric chromosomes in the human genome<sup>4-11</sup>. As a result, the T2T Consortium also provided insights into the organization and function of segdups, centromeres, epigenetic features of repeats and genome, and human genetic variation by comparative genomics and population genetics approaches<sup>4-9</sup>. These efforts significantly extend our biological understanding of human genomics and underscore the advantages of using T2T-CHM13 as a reference for genomic analyses<sup>11</sup>.

More than 200 Mbp of genomic sequences were identified as ‘non-syntenic’ regions between GRCh38 and T2T-CHM13 in the previous studies<sup>4,9</sup>, representing major of large-scale genomic differences between these two assemblies. The large-scale genomic differences are largely concentrated in complex genomic regions, which play an outstanding role in human disease as well as evolutionary adaptation<sup>10,12</sup>. For example, segdups of *Notch2NL* are associated with brain development in primate evolution, while a rare microdeletion of *Notch2NL* causes microcephaly in humans<sup>13-16</sup>. We sought to revisit this analysis using different methods in order to further refine and assess the large-scale genomic differences between GRCh38 and T2T-CHM13 for future applications (e.g., genotyping, association, and evolutionary studies).

Here, we expand on the comparison of ‘non-syntenic’ regions between T2T-CHM13 and GRCh38 in the previous studies<sup>4,9</sup>, identifying additional large-scale genomic differences between these two assemblies applying an array of additional alignment and visual validation tools. We characterize the genomic regions with at least 10 kbp genomic differences between the two assemblies into four types: insertions, deletions, inversions, and structural divergent regions (SDRs), with respect to GRCh38. We then develop an integrated website tool (SynPlotter, <https://synplotter.sjtu.edu.cn/>) to validate the discrepant regions and characterize the gene model differences in these regions. In addition, we use the 239 human genomes from the Simons Genome Diversity Project (SGDP)<sup>17</sup> to test whether these discrepant regions are likely copy number (CN) polymorphic in human populations. We also investigate the functional relationship between discrepant regions and human diseases. Finally, we systematically analyze the evolutionary history of one example of a newly identified discrepant region—the *KLRC* gene cluster—in human populations and other nonhuman primates.

## Results

### *Large-scale genomic discrepant regions*

More than 570 ‘non-syntenic’ regions (~238Mbp) have been identified between the T2T-CHM13 and GRCh38 genome assembly with a 1Mbp syntenic interval approach<sup>4,9</sup>. Here, to more completely and precisely characterize the structural types of the large-scale genomic differences between the two genome assemblies, we applied three additional alignment tools (PAV<sup>18,19</sup>, minigraph<sup>20</sup>, PBSV (<https://github.com/PacificBiosciences/pbsv>)) to expand on the non-syntenic regions originally identified by LASTZ<sup>4,9</sup>. We identified the 695 structural variants (SVs,  $\geq 10$  kbp) with three independent methods (Tables S1-S3). Next, we developed an integrated website tool (SynPlotter) that is designed to visualize and cross-validate the syntenic relationship between GRCh38 and T2T-CHM13 by integrating multiple aligners (e.g., minimap2 and numer) and publicly available visualization tools (e.g., dotplot and SafFire (unpublished, <https://mrvollger.github.io/SafFire>))<sup>21,22</sup>. Excluding the SVs in centromere and telomere regions, we validated 238 of 274 large SVs (validation rate: 86.9%) using our validation tool (Table S4).

Next, we integrated our validated large SVs with the validated ‘non-syntenic’ regions (Fig. 1a, Fig. S1, and Table S4) for a total of 590 discrepant genomic regions (~226 Mbp) between GRCh38 and T2T-CHM13 in total (Fig. 1b). Of these, 295 regions are in centromeres

(204.64 Mbp), 57 regions are in (sub)telomeres (1.23 Mbp), 162 regions are in segdups (17.86 Mbp), 18 regions are in tandem repeats (0.56 Mbp), while 58 regions occur in other parts of the genome (1.98 Mbp) (Fig. 1b).

We excluded 352 discrepant regions in centromeres and telomeres and focused on the euchromatic regions<sup>4,9,23</sup>. We refined the characterization of the large-scale discrepant regions by categorizing them into four types (including: insertions, deletions, inversions, and SDRs, with respect to GRCh38) with more precise breakpoints (Fig. 1c). There are 23 deletions (1.51 Mbp), 83 insertions (3.42 Mbp), 39 inversions (10.47 Mbp) and 26 SDRs (1.87 Mbp) in the previously reported ‘non-syntenic’ regions (total: 17.27 Mbp) (Fig. 1c). Relative to the previously reported ‘non-syntenic’ regions, here, we found 67 newly identified discrepant regions, of which the number is ~40% greater than that of the reported ‘non-syntenic’ regions (Fig. 1c). The 67 newly identified regions (total: 3.13 Mbp) include 45 deletions (1.7 Mbp), 4 insertions (0.06 Mbp), and 18 inversions (1.37 Mbp) (Fig. 1c). The number of deletions in the newly identified set is higher than in the ‘non-syntenic’ regions ( $p < 0.001$ , chi-square test).

### ***Gene model and structure differences in the CN polymorphic discrepant regions***

Of the 238 discrepant regions, 63 of them include 153 genes, such as *TBC1D3*, *AMY1*, *GPRIN2*, and *NOTCH2NL* (reported in the ‘non-syntenic’ regions)<sup>4,9</sup>. Of these, 53 protein-coding genes are in the 25 newly-identified discrepant regions, including *ZDHCC11B*, *GSTM2*, *CFHR3*, *CFHR1*, *CRI*, and *KLRC2* (Table S5). The depletion of *ZDHH11B* is found in T2T-CHM13 by a ~98 kbp deletion, with respect to GRCh38 (Fig. 2a). The read-depth genotyping from 206 Illumina short read genomes from the SGPD shows the CN polymorphism of *ZDHH11* in humans (Fig. 2a). The gene models showed that the two exons are deleted in *ZDHH11B* compared to *ZDHH11* (Fig. 2a). We also observed the depletion of *GSTM1* in T2T-CHM13 by a ~17 kbp deletion and the *GSTM* is inferred as CN polymorphic in the 206 humans (Fig. 2b). The gene models showed that a few amino acids of *GSTM1* are different from that of *GSTM2* (Fig. 2b). In addition, we observed an ~18.5 kbp deletion in T2T-CHM13, resulting in the depletion of eight exons (450 amino acids) in *CRI* (Fig. 2c). We examined the length of *CRI* gene in the 94 long-read human genome assembly from the Human Pangenome Reference Consortium (HPRC)<sup>24-26</sup> and the length of *CRI* in 79 assemblies coincides with that of T2T-CHM13. This suggests that T2T-CHM13 carries the major allele of *CRI* (allele frequency: 0.84) (Fig. S2). In addition, another ~85 kbp genomic

region, including *CFHR1* and *CFHR3*, is deleted in T2T-CHM13, with respect to GRCh38 (Fig. S3).

We assessed whether the discrepant regions are likely CN polymorphic in the human genome. We used the standard deviation (s.d.) of the CN as an index to represent the level of polymorphisms (see Methods)<sup>27</sup>. To reduce the CN estimation bias, we excluded the regions where the CN is greater than 10 in the following analyses. We observed that the mean s.d. of the CN of the 131 discrepant regions (mean=0.67) is ~5-fold greater than that of the whole genomic regions (mean=0.13, empirical p=0) (Fig. 2d and Fig. S4), as expected. We next tested whether the discrepant regions are more likely CN polymorphic than the CN variable regions (CN>2.5 and CN<10). We observed the mean s.d. of the CN of the 131 discrepant regions (mean=0.58) is ~1.2-fold greater than that of the CN variable regions (mean=0.58, empirical p=0.003) (Fig. 2e). Yet, we did not observe a significant difference between the median s.d. of the CN of the 131 discrepant regions (median=0.46) and that of the CN variable genomic regions (median=0.4, empirical p=0.07) (Fig. S4). The simulation tests imply that the discrepant regions are more likely CN polymorphic than the genome-wide average, maybe even than the CN variable regions in the human genome. These results suggest that the gene structure/model differences in the CN polymorphic discrepant regions warrant further investigation for potential disease association and functional assessment.

### ***Disease relevant loci are associated with the large-scale discrepant regions***

We integrated the reported morbid copy number variants (CNVs) and genomic disorder CNVs that associated with more than 50 disease phenotypes, including neurodevelopmental disorders, abnormality of the immune system, and others<sup>28-30</sup>. We next queried whether the discrepant regions are more likely associated with the reported disease relevant CNVs. With genome-wide permutation analysis (see Methods), we found that the discrepant regions are significantly co-localized with disease relevant CNVs (empirical p=0.003, ~1.7-fold excess) (Fig. S5). To better characterize the genes/genomic coordinates relevant to disease, we surveyed the literature and DECIPHER database for the aforementioned discrepant regions and found 27 discrepant regions associated with human diseases; 18 of them are newly identified discrepant regions (Table 1).



The genes in the 27 disease relevant discrepant regions are enriched in the neuroblastoma breakpoint family domain ( $p=3.9e-5$ ), complement and coagulation cascades ( $p=7.5e-4$ ), glutathione metabolic process ( $p=4.4e-3$ ), and antimicrobial ( $p=5.5e-5$ ) by the gene ontology (GO) enrichment analysis (Table S6). Therefore, the rare microdeletions or microduplications of these discrepant regions mainly affect the development and function of the circulatory system (urinary system disease (e.g., chromosome 1p13.3)), immune system (COVID-19 (e.g., 6p21.32, 12p13.2)), and nervous system (bipolar disorder/schizophrenia<sup>31</sup> (e.g., 10q11.22) and autism spectrum disorder (e.g., 16p12)) (Table 1). We also found some genes within the discrepant regions that are proven to be functionally well-known and pathogenic. For example, *KLRC2*, located in a newly identified discrepant regions, is involved in immune cell maturation and subtype differentiation<sup>32</sup>. The *KLRC2* protein (also: NKG2C) can bind to CD94 and HLA-E to form a functional complex<sup>33</sup>, and thus, the depletion of *KLRC2* is likely to have a significant impact on the development of severe COVID-19<sup>34</sup>. In the visual cortex, microglial CD94/ *KLRC2* is necessary for regulating the magnitude of ocular dominance plasticity during the critical period of development<sup>35</sup>. *GSTM1* (Glutathione S-Transferase Mu 1) encodes a member of a superfamily of antioxidant enzymes, which is important in kidney disease progression<sup>36</sup>. *ZDHHC11B* (Zinc Finger DHHC-Type Containing 11B) is involved in a network that promotes the proliferation of Burkitt lymphoma cells<sup>37</sup>. *CFHR3* and *CFHR1*, belonging to Complement factor H (CFH), plays an essential role in regulating the alternative pathway of the complement system<sup>38</sup>. These results suggest that the discrepant genomic regions are functionally important and need to be considered carefully with respect to genome-wide association.

### ***The diversity of KLRC2 characterized with the 94 long-read and 2,504 short-read human genomes***

We observed that *KLRC2* is deleted by a 15.4 kbp deletion variant in T2T-CHM13, with respect to GRCh38 (Fig. 3a). This discrepant region is CN polymorphic in human populations as evidenced by SGPD read-depth genotyping<sup>39-41</sup> (Fig. S6). To better characterize the diversity of the *KLRC* region, we systematically investigated the discrepant region with the 94 long-read genome assemblies from the HPRC dataset<sup>24-26</sup>. We found 1 duplication and 11 deletions of *KLRC2* in the 94 long-read genome assemblies (Table S7). We refined the breakpoints of the duplication and deletion of *KLRC2* in the T2T-CHM13 and HG002\_hap<sup>19</sup> genome assemblies at single-base pair resolution to understand the mechanisms of the structural variation.

There are four *KLRC* genes in the *KLRC* discrepant region, wherein a segdup including *KLRC2* and *KLRC3* represents is configured in a direct orientation<sup>42</sup> (Fig. 3). The configuration provides a genetic basis for microdeletions and microduplications of *KLRC2*. The syntenic relationship of the *KLRC* gene cluster between T2T-CHM13 and GRCh38 revealed that the left breakpoint of the *KLRC2* deletion is located within *AluYm1*, and the right breakpoint of the *KLRC2* deletion is 3 bp away from another *AluYm1* (Fig. 3b). We also observed a 43 bp repeat motif (tgatgcctcccaaagtctgggattataggcttgagccacca) at both breakpoints (Fig. 3b). In addition, we refined the breakpoints of the *KLRC2* duplication in the HG002\_hap1 assembly (Fig. 3c). We found that the duplication sequences (~15.4 kbp) are inserted in an *AluJb* element and the *AluJb* element is disrupted by a simple repeat insertion in GRCh38 (Fig. 3d). The breakpoints are located within poly adenine (polyA) sequences in GRCh38 (Fig. 3d).

Our analysis of the long-read HPRC haplotypes (n=94) identified three haplotypes of the *KLRC* gene cluster, including 0, 1, and 2 copies of *KLRC2*, respectively. Next, we used the SUNK (singly unique nucleotide *k*-mer) mapping and read-depth genotyping approaches<sup>39-41</sup> to infer the three haplotypes in 2,504 human genomes from the 1,000 Genome Project (1KG)<sup>43</sup>. We found that 19%, 78%, and 3% of modern humans contain 0, 1, and 2 copies of *KLRC2*, respectively (Fig. 4a, Table S8). The haplotype with a depleted *KLRC2* (“*KLRC*-hap2”) occurs more frequently in African (e.g., Esans: 25.45%) and East Asian (e.g., Kinhs: 26.24%) populations but is observed less frequently in American (e.g., Peruvians: 3.8%) and South Asian (e.g., Pakistans: 9.71%) populations (Table S9).

To study whether the depletion of *KLRC2* is a recurrent or a single-deletion event in human populations, we used the ~12.7 kbp *KLRC* gene cluster genomic regions, including both *KLRC2* and *KLRC3*, to reconstruct the phylogenetic tree of the 94 long-read human samples. The results showed that the *KLRC2* depletion haplotypes from different human groups are monophyletic (Fig. 4b), suggesting that the *KLRC*-hap2 (*KLRC2* depletion) deletion arose once in human population history.

### ***Gene expression and NK cell differentiation between two *KLRC* haplotypes***

To investigate potential functional effects of different *KLRC2* haplotypes in humans, we identified six single-nucleotide variants (SNVs) that distinguish *KLRC*-hap2 (*KLRC2*



depletion) from *KLRC*-hap1 (one copy *KLRC2*) with the 94 long-read genome assemblies. We examined the linkage disequilibrium (LD) of the *KLRC* gene cluster among 2,504 high-coverage genomes from the 1KG human population. In general, the *KLRC* gene cluster shows significant LD (LD:  $r^2 > 0.5$ ,  $D' > 0.5$ ) (Fig. 4c). In particular, the six SNVs identified in the 94 long-read genome assemblies are tightly linked (LD:  $r^2 > 0.9$ ;  $D' > 0.99$ ) (Fig. 4c). These six SNVs can, thus, be used to infer the deletion haplotype. Indeed, we find that the allele frequencies of the six SNVs of the *KLRC*-hap2 in ~135,000 humans from the gnomAD database are from 19.9% to 20.6%, which coincide with the CN frequency of *KLRC*-hap2 in humans from our above *KLRC*-haplotype inference analysis (Fig. 4c and Table S10).

Of note, the six SNVs are identical between GRCh38 and T2T-CHM13, although apparently distinguish two distinct *KLRC* haplotypes (GRCh38: *KLRC*-hap1, T2T-CHM13: *KLRC*-hap2). We investigated whether this apparent discrepancy could have resulted from a ‘mixed’ haplotype in GRCh38. In GRCh38, we note that the region was assembled by the two distinct bacterial artificial chromosome (BAC) clones (AC022075.29 and AC068775.52) from one donor (RP11) (Fig. 4c). Previous studies have shown that haplotype swaps are usually associated with the overlap boundary of the two adjacent BAC clones<sup>6</sup>. In support of this, our LD analysis within the *KLRC* locus (see Methods), shows that GRCh38 possesses combinations of alleles that are either in strong positive or strong negative LD, whereas the corresponding region of T2T-CHM13 largely exhibits alleles only in positive LD. Thus, T2T-CHM13 better reflects the haplotype structure of living human populations. LD at the *KLRC* locus extends much further than the randomly selected control locus, which exhibits multiple, shorter haplotype blocks, potentially reflecting differences in the history of recombination within the regions or a deep coalescent deletion (Fig. S7 and Fig. S8). Taken together, we consider the *KLRC* gene cluster organization in GRCh38 to represent the product of a misassembly of two different haplotypes.

Using the GTEx (Genotype-Tissue Expression) multi-tissue eQTL (expression quantitative trait loci) database (release v8, <https://gtexportal.org/>), we investigated how these six SNVs relate to *KLRC2* gene expression differences among 54 different tissues. We show that *KLRC*-hap2 SNVs correspond to reduced expression of *KLRC2* gene in 35 tissues (Fig. 4d). In particular, the brain and spleen tissues show the most significant gene expression difference between two haplotypes ( $p < 2e-5$ ). Further, we investigated the association

between these six SNVs with more than 600 phenotypes/traits (GWAS Atlas) and find that three out of six SNVs are significantly associated with immune domain function involving the NK cell differentiation ( $p < 1e-15$ )<sup>32,44</sup> (Fig. 4e). These results suggest that the depletion of *KLRC2* likely plays a role in the immune differentiation.

### ***The evolutionary history of KLRC2 and KLRC3 in primates***

Using sequence read-depth, we investigated CN of *KLRC* genes among a population of non-human primates (NHPs). Our analysis revealed that *KLRC2* and *KLRC3* are also CN polymorphic and that *KLRC2* and *KLRC3* CN in the African great apes is greater than other primates (Fig. S9). We also investigated the organization of the *KLRC* gene cluster within 16 long-read genome assemblies. The analysis confirmed *KLRC2* and *KLRC3* are CNV in NHPs with a deletion of *KLRC2* and *KLRC3* in two gibbon genome assemblies (Fig. S10) and three copies of *KLRC2* in two macaque genome assemblies (Fig. S11). We reconstructed the phylogenetic tree of *KLRC2* and *KLRC3* using ~5.5 kbp genomic region. The phylogenetic tree shows that *KLRC2* and *KLRC3* were duplicated within the common ancestor of apes and Old World monkey at ~19.8 million years ago (mya) (95% CI: 10.85-28.97 mya) (Fig. 5a and Fig. S12). In addition, we found that *KLRC2* is independently duplicated in humans and macaques (Fig. 5a).

We also examined *KLRC2* and *KLRC3* duplicated genes for evidence of natural selection during primate evolution. An analysis of diversity of the *KLRC* gene cluster based on the 94 long-read genome assemblies revealed no significant differences based on pi diversity estimates (Fig. S13). In contrast, a branch model estimate of amino acid selection, as defined by PAML and aBSREL (an adaptive branch-site REL test for episodic diversification), found evidence of selection within the *KLRC3* clade ( $p = 0.03$ , likelihood ratio test)<sup>45,46</sup> (Table S11). In particular, we identified three amino acids (R224, R227, G229) of *KLRC3* predicted to be under positive selection with greater 90% possibility by the branch-site model implemented in PAML ( $p = 0.006$ , likelihood ratio test) (Fig. 5b and Table S11).

Based on the AlphaFold2 and *KLRC1* crystal structure<sup>33,47</sup>, the protein structure of *KLRC2* and *KLRC3* are predicted to be altered by these three amino acids (Fig. 5b). *KLRC* proteins have been shown to bind CD94 and HLA-E for immune response<sup>33,47</sup>. With predicted complex protein structure of *KLRC*/HLA-E/CD94, significant differences were found between the interaction interfaces of *KLRC2*/HLA-E and *KLRC3*/HLA-E. Two hydrogen

bonds were observed between the C-terminal of KLRC3 and HLA-E: the amide nitrogen of Ile226 from KLRC3 binds the side chain of Glu175 from HLA-E; the side chain Arg227 from KLRC3 binds the carbonyl oxygen of Asp170 from HLA-E. The two hydrogen bonds may stabilize the flexible loop of KLRC3 (Fig. 5c). However, no obvious interactions were found between the C-terminal of KLRC2 and HLA-E (Fig. 5c). Our findings implicate differential binding affinity at the two interaction interfaces potentially important for functional differentiation of KLRC2 and KLRC3 in humans.

## Discussion

The first complete genome assembly (T2T-CHM13) represents a new and important genomics resource<sup>4,11,48,49</sup>. Here, we more systematically investigate the large-scale genomic differences between T2T-CHM13 and the current reference genome assembly (GRCh38). We show that the discrepant regions are among the most structurally complex and may introduce reference biases in human genetics (e.g., genotype-phenotype association study) and evolutionary genomics (e.g., gene family evolution investigation). Therefore, understanding the discrepant regions between the two crucial reference genome assemblies will provide a crucial resource for further genomic and functional studies. In this study, we systematically characterized the large SVs between the two human genome assemblies and found 67 newly identified discrepant regions. In addition, we developed an integrated website tool (SynPlotter) to visualize and validate 246 discrepant regions. The newly identified regions include gene-model differences (e.g., *ZDCHH11B*, *GSTM1*, *CFHR3*, *CFHR1*, *CR1*, and *KLRC2*) and the SGDP read-depth genotyping data show that the discrepant regions are more likely CN polymorphic. In addition, the discrepant regions are often related to human diseases. Finally, we provided a novel example to illustrate the biological importance of discrepant regions by analyzing the *KLRC* gene cluster with population genetics and evolutionary genomic approaches.

Previous studies used a 1Mbp syntenic intervals to identify the sequence difference between the two genome assemblies<sup>4,9</sup>. Here, we used three different methods to identify the SVs ( $\geq 10$  kbp) with reciprocal alignment between GRCh38 and T2T-CHM13 to identify precise breakpoints and structural types of the large-scale discrepant regions (Fig. S1). We additionally identified the 67 discrepant regions and ~70% of them belongs to deletions, because deletion variants are likely chained in a large synteny by the LASTZ<sup>4,9</sup>. Notably, comparing with the recent inversion dataset in humans<sup>50</sup>, we found that one inversion in our

dataset was not reported in the dataset<sup>50</sup>. The genomic region of the inversion contains a gap in GRCh38 (Fig. S14). This shows the SV discovery would be affected by the reference bias and T2T-CHM13 is useful to identify large-scale SVs. In addition, we developed an integrated visualization tool to validate the discrepant regions. This website tool is user-friendly and publicly available to compare syntenic regions between GRCh38 and T2T-CHM13.

The discrepant regions between these two assemblies have been regarded as CN polymorphic genomic regions in previous studies<sup>4,6,9</sup>. We performed CN analysis to provide clear evidence to support that the discrepant regions are likely more polymorphic than the genome-wide average, even than the CN variable regions. In addition, to our knowledge, we surveyed the relevance between the discrepant regions and the reported medical relevant loci in greater detail. We find that rare microdeletions and microduplications of 27 discrepant regions are potentially related to neurodevelopmental diseases and others with supported evidence<sup>28-30</sup>. Loss of function of *CRI* is associated with the Alzheimer's disease<sup>51,52</sup> and T2T-CHM13 carries a major allele of *CRI*. Yet, GRCh38 carries a minor allele, where eight exons encoding tandem repeat protein domain in *CRI* are inserted with respect to T2T-CHM13. In addition, *ZDHHC11* (Zinc Finger DHHC-Type Containing 11) and *ZDHHC11B* are involved in innate immune or anti-virus response by enabling signaling adaptor activity. The CNV of *ZDHHC11* and *ZDHHC11B* are associated with hepatoblastoma<sup>53</sup> and primary open-angle glaucoma<sup>54</sup>. The *GSTM1* (glutathione S-transferase mu) locus is also a polymorphic locus associated with cancers, metabolism, and hepatic cirrhosis<sup>55</sup>. Thus, our study provides a fundamental resource for functional assessments to examine functional differentiation between/among polymorphic loci in humans. Importantly, it is still unclear whether the reference bias has effects on the reported disease association study of these discrepant regions. If so, the excess of the co-localization between the discrepant regions and disease relevant CNV needs to be re-assessed.

We comprehensively compare the *KLRC* gene cluster in humans and NHPs. Firstly, we precisely characterize the breakpoints of duplication and deletion of *KLRC2*. The breakpoints on single-base-pair resolution could facilitate the molecular probe development to genotype the CN of *KLRC2* in the future. The duplication and deletion mechanism of *KLRC2* are associated with the *Alu* elements and simple repeats in the human genome. Notably, our *KLRC* haplotype inference and phylogenetic tree analyses show that the origin of *KLRC2* and

*KLRC3* is duplicated from the common ancestor of the apes and Old World monkey. The human population genetic analyses reveal that *KLRC*-hap2 (*KLRC2* depletion) is caused by a single deletion event in humans. Africans and Asians have a higher frequency of *KLRC*-hap2 but we did not observe significant pi diversity change in the *KLRC* gene cluster in humans. These results would suggest that the distribution of the *KLRC* haplotypes may simply be the result of genetic drift in human evolution. Yet, we identified the six SNVs to distinguish *KLRC*-hap2 (*KLRC2* depletion) from *KLRC*-hap1. The eQTL and GWAS analyses show the gene expression and immune functional differentiation between the two *KLRC* haplotypes, and previous functional experiments shows the *KLRC2* haplotypes have different roles in synaptic pruning<sup>35</sup>. Additional experiments are required to determine if loss of *KLRC2* is the result of genetic drift or subject to other models of selection (e.g., balancing selection).

Our tests of selection implicate three amino acids of *KLRC3* as potentially subject to positive selection during primate evolution. Predicted protein structures further suggest structural differences (KLRCX/HLA-E/CD94) between *KLRC2* (KLRC2/HLA-E/CD94) and *KLRC3* (KLRC3/HLA-E/CD94). It is possible that *KLRC3* has acquired distinct functional properties from *KLRC2* as a result of natural selection.

We also show that the *KLRC* gene cluster region was misassembled in GRCh38, likely because the region was assembled by two BAC clones from two distinct *KLRC* haplotypes. If we used the six SNVs to infer the *KLRC* haplotype, GRCh38 would carry *KLRC*-hap2 (*KLRC2* deletion). Yet, the GRCh38 shows *KLRC*-hap1 at present. As a result, association studies of *KLRC* genes and their interpretation would be potentially confounded.

Altogether, our study provides a more comprehensive and detailed assessment of the structure and function of the large-scale discrepant genomic regions between GRCh38 and T2T-CHM13. We believe the results of this work not only contribute to our biological understanding of these diverse regions but will benefit future studies by helping to eliminate reference biases. We should stress that our study focused solely on the large-scale discrepant regions between two ‘completed’ genome assemblies and, as such, represents a limited survey of the true extent of human genome structural variation<sup>18,19</sup>. It is anticipated that the T2T Consortium will generate more complete genome assemblies from a diversity of human samples and non-human primates. These will help us to fully understand the extent of

complex/discrepant regions in humans<sup>4-11,25,26</sup> and their biological impact using reference-free approaches.

## Methods

### Data in this study

We downloaded 94 long-read human genome assemblies from the HPRC phase 1 project (<https://humanpangenome.org/>)<sup>24-26</sup>. We download the Illumina data of the 2,504 high-coverage short-read from the 1KG human population dataset<sup>43</sup>. For the reconstruction of the phylogeny of the *KLRC* gene cluster, we locally assembled the *KLRC* gene cluster region from the published HiFi reads of chimpanzee, bonobo, gorilla, orangutan, gibbon, macaque, owl monkey, and marmoset. In addition, the ‘non-syntenic’ region, centromere, and gene annotation files were downloaded from the UCSC Genome Browser directly (<https://genome.ucsc.edu>). The (sub)telomere regions are defined as a 500 kbp region away from the start or end of chromosome in this study.

### Discrepant region characterization and validation

We used a reciprocal alignment approach to systematically characterize the SVs. In detail, We used GRCh38 as the reference genome and T2T-CHM13 as the query to run PAV (v2.0.0)<sup>18,19</sup>, PBSV (PBSV, <https://github.com/PacificBiosciences/pbsv>, v2.8.0), and minigraph<sup>20</sup> to characterize SVs. We also used T2T-CHM13 as the reference and GRCh38 as the query to run PAV, PBSV, and minigraph to characterize SVs. Then, we LiftOver the coordinates from GRCh38 to T2T-CHM13 and then merged these calls with bedtools (v.0.29.0)<sup>56</sup>. In the PBSV analysis, we used the PBSIM2 tool<sup>57</sup> to simulate HiFi reads from GRCh38 and used these simulated reads against T2T-CHM13.

We developed a custom script to automatically screenshot the syntenic plots from SafFire (<https://mrvollger.github.io/SafFire>). In addition, we integrated minimap2 (v2.24)<sup>21</sup> or mummer4 aligner<sup>22</sup> to generate syntenic PAF files. We next applied the dotplot implement in mummer4<sup>22</sup> to generate dot plots. Then, we implement the above scripts into a website tool (SynPlotter) to visualize the syntenic relationship between two coordinates. With our website tool, the syntenic plot and dot plot can be generated, and the basic genomic difference statistics could be calculated. The gene and repeat (e.g., segdups) annotations are also shown.



Lastly, we used our website tool to validate the large SVs ( $\geq 10$  kbp) generated by the three above callers.

### **Gene, structural type, and repetitive sequence annotation for discrepant regions**

We used bedtools (v2.29.0)<sup>56</sup> to compare the discrepant regions between ‘non-syntenic’ regions and our large SVs ( $\geq 10$  kbp). We used the ggplot2<sup>58</sup> and karyoploteR packages<sup>59</sup> to plot the chromosome ideogram. Next, we characterized the discrepant regions into four types (insertions, deletions, inversions, and SDRs) with eyes and manually refined the breakpoints of these SVs.

We also used the gene, repetitive sequence annotation files from the UCSC Genome Browser to annotate these discrepant regions with bedtools (v2.29.0)<sup>56</sup>. In this study, the centromere regions include pericentric regions, but not the centromeric transition regions in all analysis. Notably, the previous study reported sequence difference between GRCh38 and CHM13, but our study reported the location of the discrepant regions. In the Fig. 1, we counted the number and the length of the discrepant regions as centromere (CEN), segdup (SD), telomere (TEL), and tandem repeat (TRF) based on the location of them, not the absolute length of sequences belonging to each type.

We used a hierarchy approach to define SV-location: (1). If a given SV located in CEN, we counted it as CEN; (2). If a given SV located in TEL (500kbp from head or tail), we counted it as TEL; (3). If a given SV located in segdup and includes at least 20% or 2kbp segdup sequence, we counted it as SD; (4). If a given SV located in TRF and includes at least 10% or 1kbp sequence, we counted it as TRF; (5). If a given SV does not belong to any type of above, we counted it as others. For example, if a given SV located in the centromere regions, we counted it as CEN type no matter whether it contains segdup sequence or not.

### **Structural polymorphism enrichment analysis**

To test whether the discrepant regions are more likely polymorphic, we downloaded the SGDP CN table<sup>17</sup> from the UCSC Genome Browser (<https://genome.ucsc.edu/>). Here, to reduce the bias from the high CN (average CN estimation from SGDP  $\geq 10$ ,  $n=53$ ), we only used the 131 discrepant regions (CN  $< 10$ ) belonging to insertions, deletions, and SDRs to calculate the standard deviations (s.d.) of the CN. First, we used bedtools (v.2.29.0)<sup>56</sup> to intersect the 131 discrepant regions with the SGDP CN table and then calculated the s.d. of

the CN of each intersected region. Then, we calculated the mean and median values of estimated s.d. (mean=0.735, median=0.439). These are our observed s.d. values of the 131 discrepant regions.

For null distribution, we did two experiments in this study. (1). Simulate the distribution of mean s.d. of the whole genome. We used bedtools (v2.29.0)<sup>56</sup> to randomly shuffle the corresponding number of coordinates (n=131) in the genomic regions where there are no centromeres, telomeres, and CN<10. We intersected them with the SGDP CN table, and calculated the s.d. of the CN for each intersected region. Then, we calculated the mean/median s.d. value. (2). Simulate the distribution of mean s.d. of the CNV regions (CN>2.5 and CN<10). We used bedtools (v2.29.0) to randomly shuffle the corresponding number of coordinates (n=126) in the genomic regions where there are no centromeres, telomeres, CN >2.5, and CN<10. We intersected them with the SGDP CN table, and calculated the s.d. of the CN for each intersected region. Then, we calculated the mean/median s.d. value.

We repeated this 1000 times for each experiment and calculated the empirical p-value of our observed mean s.d. value (permutation test). In addition, we also estimated the observed mean and median s.d. values of the SGDP CN table for the two experiments. (1). We calculated the mean and median s.d. value of the regions (CN <10) based the SGDP CN table (mean = 0.13, median =0.079). (2). We calculated the mean and median s.d. value of the regions (CN <10 and CN >2.5) based on the SGDP CN table (mean = 0.58, median =0.463).

### **Disease relevant CNV enrichment and survey**

We downloaded the CNV coordinates from the morbid and the cross-disorder dosage sensitivity maps<sup>28-30</sup> and LiftOver them to T2T-CHM13. We next used bedtools to intersect our discrepant regions with the integrated coordinates and found 33 discrepant regions are co-localized with the disease relevant CNVs. We also used a permutation test to shuffle the discrepant regions in T2T-CHM13, excluding centromeres and telomeres, and calculated how many discrepant regions could be co-localized with the disease relevant CNVs. We repeated this process 1,000 times and plotted the distribution of the number of co-localized regions (mean N=19.9) with ggplot2 in R.

We also manually curated the coordinates with gene annotations from the literatures and DECIPHER database by hands. To better represent the disease relevant discrepant regions (Table S5), we only listed the regions with well-qualified evidence/literature/case-report to support as disease relevant in Table 1.

### **Genomic syntenic comparison analysis**

In this study, we found 67 newly identified discrepant regions, of these, the 25 regions contain 38 genes. Thus, we used minimiro (commit 18271297374ae6a679521a7ce3f5bb6c0cf8d261) to compare the genomic syntenic relationships between GRCh38 and T2T-CHM13 in these regions.

Then, we used the RefSeq annotation from GRCh38 and CAT/RefSeq annotation from T2T-CHM13 to extract the protein sequences of the genes. The mafft program (v7.4.3)<sup>60</sup> was used to align the amino acid to check the amino acid difference among the homologous genes. The schematic plots were generated by ggplot2 and AliView (v1.26)<sup>61</sup>.

### ***KLRC2* haplotype characterization**

We extracted the genomic regions containing the *KLRC* gene cluster region from GRCh38 (chr12:10359648-10470462). Then, we used minimap2 (v2.24) to map the region to 94 long-read human genome assemblies and other NHP long-read genome assemblies. Finally, we used the minimiro to generate the syntenic plots between GRCh38 and other human and NHP samples. We found three distinguished haplotypes of the *KLRC* gene cluster based on the CN variation of *KLRC2*.

### ***KLRC* haplotype inference from 1KG population dataset**

We used the previously reported read-mapping approach (SUNK-WSSD)<sup>39,40</sup> to genotype the CN of *KLRC2* and *KLRC3* in the 1KG population dataset (2,504 high-coverage Illumina genomes). The mean CN of *KLRC3* of 2,504 humans is ~1.8 (s.d.: 0.14) (Fig. S4), while the mean CN of *KLRC2* of 2,504 humans is ~1.9 (s.d.: 0.6). The CN of *KLRC3* is clustered together in different human groups, suggesting there is no CN variation of *KLRC3* in humans (Fig. S15). However, the CN of *KLRC2* is clustered into three groups in different human groups, suggesting there is CN variation of *KLRC2* in humans.

If the CN of *KLRC2* is greater than 2.5 (mean + 1 s.d.), we inferred the *KLRC2* CN as 3. If the CN of *KLRC2* is less than 1.3 (mean – 1 s.d.), we inferred the *KLRC2* CN as 1. Then, we used the maps, ggplot2, and scatterpie packages (<https://www.rdocumentation.org/packages/scatterpie/versions/0.1.8>) in the R to plot the world map of the *KLRC* haplotype map.

## Phylogeny reconstruction and time-calibration tree reconstruction

We used minimap2 (v2.4) to determine the syntenic regions in human and NHP genome assemblies. We also used samtools (v1.9) to extract the corresponding regions. Then, we used mafft (v7.453) to align the genomic sequence with default parameters and used it as input for IQTREE (v1.6.11) to build the maximum likelihood phylogenetic trees<sup>21,62-64</sup>.

To determine the time of duplication, we used BEAST2 (v2.6.2)<sup>65</sup> to date the phylogeny. The prior calibrated times were used from the previous studies. Here, we used the log-normal and the real mean model to set the prior calibrate time, including pan-lineage split time (~1.45 mya), owl monkey and marmoset split time (~24.5 mya), monkey and ape split time (~54 mya) in this study. The MCMC chains were run 30,000,000 steps and 3,000,000 steps were set for burnin running. Finally, we used the tracer (v1.7.1) to examine whether the chain was convergent. Indeed, each ESS value of each parameter was over 200 in our study and these results suggested the MCMC chain was converged. FigTree (v1.4.4) and iTol (<https://itol.embl.de/>) were used to plot the trees.

## Selection test with PAML and aBSREL

We downloaded the human and NHP coding sequences (CDS) and protein sequences for *KLRC2* and *KLRC3* from the UCSC Genome Browser. We used mafft (v7.4.3) to align the protein sequences and used translatorex\_vLocal.pl to align the CDS based on the aligned protein sequences. All protocols are based on the previously reported tool (TREEasy)<sup>62-64</sup>.

We examined the pi diversity of the *KLRC* gene cluster regions in humans with 94 long-read genome genotyping data. Then, we ran a preliminary selection test on aBSREL<sup>46</sup> (<https://www.datamonkey.org/analyses>) and the aBSREL tool showed the selection signals on the *KLRC3* clade. We also ran the branch model in PAML (v4.9)<sup>45,66</sup> and the model showed the selection signals on the *KLRC3* clade too (p=0.03). The branch-site model in PAML (v4.9) shows three amino acids under selection with a probability greater than 0.9 in

the clade of KLRC3 in the Bayes Empirical Bayes (BEB) analysis ( $p=0.006$ ). The p-values were calculated by the likelihood ratio test in R.

### **Protein structure analysis of KLRC2 and KLRC3**

We predicted the structural model of KLRC2 (residues 118-231) and KLRC3 (residues 118-240) using AlphaFold2<sup>67</sup> and KLRC1 crystal structure<sup>33,47</sup> with predicted local distance difference test (pLDDT) values as 92.74 and 80.55, respectively, suggesting that they are accurate enough for the further analysis. AlphaFold2<sup>67</sup> was used to predict the structures of KLRC2/HLA-E/CD94 and KLRC3/HLA-E/CD94 complexes. The pLDDT values of the two complexes are 80.39 and 76.05, respectively, which are of high confidences and are accurate enough for the interaction analysis. Protein structure and interaction analyses were performed on PyMol (v2.4.1, <https://pymol.org/>). Structure alignment shows obvious differences between KLRC2 and KLRC3. Met223 of KLRC2 and Arg224, Arg227, and Gly229 of KLRC3 are located at the surface loops which connecting two  $\beta$ -strands, and the loop of KLRC3 has a longer conformation. His226 and Lys228 of KLRC2 may contribute to the following  $\beta$ -strand, which is longer than that of KLRC3.

### **eQTL analysis and GWAS ATLAS analysis**

We firstly aligned the *KLRC* gene cluster of 94 long-read human genome assemblies and we used our custom script to find the SNVs that are different between *KLRC*-hap2 and *KLRC*-hap1. Then, we investigated LD among Lewontin's  $D'$  and  $R^2$  implemented in LDBlockShow (v1.40)<sup>68</sup> and PLINK (v1.90b6.21)<sup>69</sup> with 2,504 high-coverage genotyping data from the 1KG dataset. The SNVs with minor allele frequencies  $>10\%$  were used for this analysis. We also calculated the allele frequency of the six distinguished SNVs in the gnomAD dataset (v.3.1.2, <https://gnomad.broadinstitute.org/>)<sup>70</sup>. The LD heatmaps were generated by LDBlockShow or R.

We also used PLINK (v1.90b6.21)<sup>69</sup> to compute LD (measured with  $D'$ ) among all SNV pairs. For each SNV pair, we then compared the reference alleles to the combinations of alleles that were determined to be “in phase” (i.e., observed together on haplotypes more often than expected under linkage equilibrium). For cases where the reference genome carried alleles that were in phase,  $D'$  was retained as a positive value, whereas for cases where the reference genome carried alleles that were out of phase,  $D'$  was multiplied by -1 to

indicate that the alleles are in repulsion. The same procedure was repeated for the corresponding region of T2T-CHM13 as determined by LiftOver, using genotype data produced by Aganezov et al<sup>6</sup>. For comparison, we also performed the same analysis for a randomly selected “control” region of the same length (83.7 kbp) for both GRCh38 and T2T-CHM13.

We downloaded the eQTL multi-tissue data from GTEx (release v8, <https://gtexportal.org/>) and we extracted the gene expression associated with the six SNVs in different tissues. The data showed that the six SNVs are only associated with *KLRC2* gene expression, as we expected. Then, we used the ggplot2 package in R to plot the normalized effect size and p-values of gene expression difference by the six SNVs. In addition, we also investigated whether any locus is related to the reported genome-wide association study (GWAS). Then, we download the phenome-wide association studies (PheWAS) data from GWAS ATLAS<sup>44</sup>. The data showed that three of the six SNVs are significantly associated with NK cells (NKearly: %335+314- and NKeff: %314-R7-)<sup>32</sup>.

## Data availability

The website tool SynPlotter is available at <https://synplotter.sjtu.edu.cn/>. The syntenic comparison of discrepant regions is deposited in GitHub ([https://github.com/YafeiMaoLab/discrepant\\_region](https://github.com/YafeiMaoLab/discrepant_region)).

## Acknowledgements

We thank T. Brown for assistance in editing this manuscript. We thank P. Hsieh and H. Cheng for their comments and scripts. We thank the HPRC for providing the 94 long-read human genome assemblies. The computations in this paper were run on the Siyuan-1 and  $\pi$  2.0 cluster supported by the Center for High Performance Computing at Shanghai Jiao Tong University.

## Funding

This work was supported by Shanghai Pujiang Program (22PJ1407300) and Shanghai Jiao Tong University 2030 Program (C-type) to Y.M.; and by National Natural Science Foundation of China (82001372) to X.Y. This work was supported by Opening research fund from Shanghai Key Laboratory of Stomatology, Shanghai Ninth People’s Hospital, College of Stomatology, Shanghai Jiao Tong University School of Medicine (Grant No. 2022SKLS-



KFKT007) to Y.M. This work was supported, in part, by US National Institutes of Health (NIH) grant# HG002385 to E.E.E. Additional funding included: NIGMS grants: K99GM147352 (to G.A.L.) and National Natural Science Foundation of China grant: 32000812 (to J.S.). E.E.E. is an investigator of the Howard Hughes Medical Institute.

## Contributions

Y.M. conceived the project. X.W., X.Y. and Y.M. finalized the manuscript. X.W., and Y.Z. performed the SVs analysis. X.W., Y.Z., and X.Y. performed the analysis of large discrepant regions. X.W., X.Y., and Y.M. performed the *KLRC* gene cluster analysis. M.X., Q.L., and Y.M. performed the *KLRC* protein structure analysis. Y.M., D.J.T., R.C.M., and M.C.S. performed the *KLRC2* haplotype swap analysis in GRCh38. S.Z. built the ‘SynPlotter’ website. Y.M., X.Y., X.W., M.X., M.R.V., N.C., W.T.H., G.A.L., D.M., J.S., R.C.M., M.C.S., W.L., Q.L., and E.E.E. contributed to interpret results and edited the draft manuscript. All authors read, edited and approved the manuscript.

## Ethics declarations

Conflict of interest

E.E.E. is a scientific advisory board (SAB) member of Variant Bio. N.C. is a full-time employee of Exai Bio. The other authors declare that they have no conflict of interest.

## Reference

- 1 Lander, E. S. *et al.* Initial sequencing and analysis of the human genome. *Nature* **409** (2001).
- 2 Venter, J. C. *et al.* The sequence of the human genome. *Science* **291**, 1304-1351 (2001).
- 3 Navarro Gonzalez, J. *et al.* The UCSC genome browser database: 2021 update. *Nucleic Acids Res.* **49**, D1046-D1057 (2021).
- 4 Nurk, S. *et al.* The complete sequence of a human genome. *Science* **376**, 44-53 (2022).
- 5 Altemose, N. *et al.* Complete genomic and epigenetic maps of human centromeres. *Science* **376**, eabl4178 (2022).
- 6 Aganezov, S. *et al.* A complete reference genome improves analysis of human genetic variation. *Science* **376**, eabl3533 (2022).
- 7 Gershman, A. *et al.* Epigenetic patterns in a complete human genome. *Science* **376**, eabj5089 (2022).
- 8 Hoyt, S. J. *et al.* From telomere to telomere: The transcriptional and epigenetic state of human repeat elements. *Science* **376**, eabk3112 (2022).
- 9 Vollger, M. R. *et al.* Segmental duplications and their variation in a complete human genome. *Science* **376**, eabj6965 (2022).

- 10 Logsdon, G. A., Vollger, M. R. & Eichler, E. E. Long-read human genome sequencing and its applications. *Nat. Rev. Genet.* **21**, 597-614 (2020).
- 11 Mao, Y. & Zhang, G. A complete, telomere-to-telomere human genome sequence presents new opportunities for evolutionary genomics. *Nat. Methods* **19**, 635-638 (2022).
- 12 Eichler, E. E. Genetic variation, comparative genomics, and the diagnosis of disease. *N. Engl. J. Med.* **381**, 64-74 (2019).
- 13 Fiddes, I. T. *et al.* Human-specific NOTCH2NL genes affect notch signaling and cortical neurogenesis. *Cell* **173**, 1356-1369. e1322 (2018).
- 14 Suzuki, I. K. *et al.* Human-specific NOTCH2NL genes expand cortical neurogenesis through delta/notch regulation. *Cell* **173**, 1370-1384. e1316 (2018).
- 15 Ishiura, H. *et al.* Noncoding CGG repeat expansions in neuronal intranuclear inclusion disease, oculopharyngodistal myopathy and an overlapping disease. *Nat. Genet.* **51**, 1222-1232 (2019).
- 16 Sone, J. *et al.* Long-read sequencing identifies GGC repeat expansions in NOTCH2NLC associated with neuronal intranuclear inclusion disease. *Nat. Genet.* **51**, 1215-1221 (2019).
- 17 Watkins, W. S. *et al.* The Simons Genome Diversity Project: a global analysis of mobile element diversity. *Genome Biol. Evol.* **12**, 779-794 (2020).
- 18 Ebert, P. *et al.* Haplotype-resolved diverse human genomes and integrated analysis of structural variation. *Science* **372**, eabf7117 (2021).
- 19 Porubsky, D. *et al.* Gaps and complex structurally variant loci in phased genome assemblies. *bioRxiv* (2022).
- 20 Li, H., Feng, X. & Chu, C. The design and construction of reference pangenome graphs with minigraph. *Genome Biol.* **21**, 1-19 (2020).
- 21 Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**, 3094-3100 (2018).
- 22 Marçais, G. *et al.* MUMmer4: A fast and versatile genome alignment system. *PLoS Comput. Biol.* **14**, e1005944 (2018).
- 23 Logsdon, G. A. *et al.* The structure, function and evolution of a complete human chromosome 8. *Nature* **593**, 101-107 (2021).
- 24 Cann, H. M. *et al.* A human genome diversity cell line panel. *Science* **296**, 261-262 (2002).
- 25 Jarvis, E. D. *et al.* Semi-automated assembly of high-quality diploid human reference genomes. *Nature* (2022).
- 26 Wang, T. *et al.* The Human Pangenome Project: a global resource to map genomic diversity. *Nature* **604**, 437-446 (2022).
- 27 Dennis, M. Y. *et al.* The evolution and population diversity of human-specific segmental duplications. *Nat. Ecol. Evol.* **1**, 1-10 (2017).
- 28 Cooper, G. M. *et al.* A copy number variation morbidity map of developmental delay. *Nat. Genet.* **43**, 838-846 (2011).
- 29 Buttermore, E. *et al.* Neurodevelopmental copy-number variants: A roadmap to improving outcomes by uniting patient advocates, researchers, and clinicians for collective impact. *Am. J. Hum. Genet.* **109**, 1353-1365 (2022).
- 30 Collins, R. L. *et al.* A cross-disorder dosage sensitivity map of the human genome. *Cell* **185**, 3041-3055.e3025 (2022).
- 31 Chen, J. *et al.* A pilot study on commonality and specificity of copy number variants in schizophrenia and bipolar disorder. *Transl. Psychiatry* **6**, e824-e824 (2016).
- 32 Roederer, M. *et al.* The genetic architecture of the human immune system: a bioresource for autoimmunity and disease pathogenesis. *Cell* **161**, 387-403 (2015).

- 33 Kaiser, B. K., Pizarro, J. C., Kerns, J. & Strong, R. K. Structural basis for NKG2A/CD94 recognition of HLA-E. *Proc. Natl. Acad. Sci. USA* **105**, 6696-6701 (2008).
- 34 Vietzen, H. *et al.* Deletion of the NKG2C receptor encoding KLRC2 gene and HLA-E variants are risk factors for severe COVID-19. *Genet. Med.* **23**, 963-967 (2021).
- 35 Marin, I. A. *et al.* The nonclassical MHC class I Qa-1 expressed in layer 6 neurons regulates activity-dependent plasticity via microglial CD94/NKG2 in the cortex. *Proc. Natl. Acad. Sci. USA* **119**, e2203965119 (2022).
- 36 Gigliotti, J. C. *et al.* GSTM1 Deletion Exaggerates Kidney Injury in Experimental Mouse Models and Confers the Protective Effect of Cruciferous Vegetables in Mice and Humans. *J. Am. Soc. Nephrol.* **31**, 102-116 (2020).
- 37 Dzikiewicz-Krawczyk, A. *et al.* ZDHHC11 and ZDHHC11B are critical novel components of the oncogenic MYC-miR-150-MYB network in Burkitt lymphoma. *Leukemia* **31**, 1470-1473 (2017).
- 38 Tschernoster, N. *et al.* Unraveling Structural Rearrangements of the CFH Gene Cluster in Atypical Hemolytic Uremic Syndrome Patients Using Molecular Combing and Long-Fragment Targeted Sequencing. *J. Mol. Diagn.* **24**, 619-631 (2022).
- 39 Mao, Y. *et al.* A high-quality bonobo genome refines the analysis of hominid evolution. *Nature* **594**, 77-81 (2021).
- 40 Sudmant, P. H. *et al.* Diversity of human copy number variation and multicopy genes. *Science* **330**, 641-646 (2010).
- 41 Sudmant, P. H. *et al.* Evolution and diversity of copy number variation in the great ape lineage. *Genome Res.* **23**, 1373-1382 (2013).
- 42 Jiang, Z., Hubley, R., Smit, A. & Eichler, E. E. DupMasker: a tool for annotating primate segmental duplications. *Genome Res.* **18**, 1362-1368 (2008).
- 43 McVean, G. A. *et al.* An integrated map of genetic variation from 1,092 human genomes. *Nature* **491**, 56-65 (2012).
- 44 Watanabe, K. *et al.* A global overview of pleiotropy and genetic architecture in complex traits. *Nat. Genet.* **51**, 1339-1348 (2019).
- 45 Yang, Z. PAML 4: phylogenetic analysis by maximum likelihood. *Mol. Biol. Evol.* **24**, 1586-1591 (2007).
- 46 Smith, M. D. *et al.* Less is more: an adaptive branch-site random effects model for efficient detection of episodic diversifying selection. *Mol. Biol. Evol.* **32**, 1342-1353 (2015).
- 47 Petrie, E. J. *et al.* CD94-NKG2A recognition of human leukocyte antigen (HLA)-E bound to an HLA class I leader sequence. *J. Exp. Med.* **205**, 725-735, (2008).
- 48 Halldorsson, Bjarni V., *et al.* "The sequences of 150,119 genomes in the UK Biobank." *Nature* 607.7920, 732-740 (2022).
- 49 Noyes, Michelle D., *et al.* "Familial long-read sequencing increases yield of de novo mutations." *The American Journal of Human Genetics* 109.4, 631-646 (2022).
- 50 Porubsky, D. *et al.* Recurrent inversion polymorphisms in humans associate with genetic instability and genomic disorders. *Cell* **185**, 1986-2005. e1926 (2022).
- 51 Brouwers, N. *et al.* Alzheimer risk associated with a copy number variation in the complement receptor 1 increasing C3b/C4b binding sites. *Mol. Psychiatry* **17**, 223-233 (2012).
- 52 Lambert, J.-C. *et al.* Genome-wide association study identifies variants at CLU and CR1 associated with Alzheimer's disease. *Nat. Genet.* **41**, 1094-1099 (2009).
- 53 Wu, J.-F. *et al.* Copy-number variations in hepatoblastoma associate with unique clinical features. *Hepatol. Int.* **7**, 208-214 (2013).

- 54 Lo Faro, V., Ten Brink, J. B., Snieder, H., Jansonius, N. M. & Bergen, A. A. Genome-wide CNV investigation suggests a role for cadherin, Wnt, and p53 pathways in primary open-angle glaucoma. *BMC genomics* **22**, 1-20 (2021).
- 55 Gu, Y. *et al.* The influence of polymorphic GSTM1 gene on the increased susceptibility of non-viral hepatic cirrhosis: evidence from observational studies. *Eur. J. Med. Res.* **23**, 1-9 (2018).
- 56 Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* **26**, 841-842 (2010).
- 57 Ono, Y., Asai, K. & Hamada, M. PBSIM2: a simulator for long-read sequencers with a novel generative model of quality scores. *Bioinformatics* **37**, 589-595 (2021).
- 58 Ginestet, C. ggplot2: elegant graphics for data analysis. *J. R. Stat. Soc. Ser. A Stat. Soc.* **174**, 245-245 (2011).
- 59 Gel, B. & Serra, E. karyoploteR: an R/Bioconductor package to plot customizable genomes displaying arbitrary data. *Bioinformatics* **33**, 3088-3090 (2017).
- 60 Katoh, K., Misawa, K., Kuma, K. i. & Miyata, T. MAFFT: a novel method for rapid multiple sequence alignment based on fast Fourier transform. *Nucleic Acids Res.* **30**, 3059-3066 (2002).
- 61 Larsson, A. AliView: a fast and lightweight alignment viewer and editor for large datasets. *Bioinformatics* **30**, 3276-3278 (2014).
- 62 Li, H. *et al.* The sequence alignment/map format and SAMtools. *Bioinformatics* **25**, 2078-2079 (2009).
- 63 Nguyen, L.-T., Schmidt, H. A., Von Haeseler, A. & Minh, B. Q. IQ-TREE: a fast and effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol. Biol. Evol.* **32**, 268-274 (2015).
- 64 Mao, Y., Hou, S., Shi, J. & Economo, E. P. TREEasy: An automated workflow to infer gene trees, species trees, and phylogenetic networks from multilocus data. *Mol. Ecol. Resour.* **20** (2020).
- 65 Bouckaert, R. *et al.* BEAST 2: a software platform for Bayesian evolutionary analysis. *PLoS Comput. Biol.* **10**, e1003537 (2014).
- 66 Yang, Z. PAML: a program package for phylogenetic analysis by maximum likelihood. *Comput. Appl. Biosci.* **13**, 555-556 (1997).
- 67 Jumper, J. *et al.* Highly accurate protein structure prediction with AlphaFold. *Nature* **596**, 583-589 (2021).
- 68 Dong, S.-S. *et al.* LDBlockShow: a fast and convenient tool for visualizing linkage disequilibrium and haplotype blocks based on variant call format files. *Brief. Bioinform.* **22** (2021).
- 69 Purcell, S. *et al.* PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am. J. Hum. Genet.* **81**, 559-575 (2007).
- 70 Collins, R. L. *et al.* A structural variation reference for medical and population genetics. *Nature* **581**, 444-451 (2020).
- 71 Byman, E. *et al.* Alpha-amylase 1A copy number variants and the association with memory performance and Alzheimer's dementia. *Alzheimers Res. Ther.* **12**, 158 (2020).
- 72 Davis, J. M., Heft, I., Scherer, S. W. & Sikela, J. M. A Third Linear Association Between Olduvai (DUF1220) Copy Number and Severity of the Classic Symptoms of Inherited Autism. *Am. J. Psychiatry* **176**, 643-650 (2019).
- 73 Wu, L. *et al.* Copy number variations of HLA-DRB5 is associated with systemic lupus erythematosus risk in Chinese Han population. *Acta Biochim. Biophys. Sin. (Shanghai)* **46**, 155-160 (2014).

832 74 Wu, Z. *et al.* Copy number variation of the Lipoprotein(a) (LPA) gene is associated  
833 with coronary artery disease in a southern Han Chinese population. *Int. J. Clin. Exp.*  
834 *Med.* **7**, 3669-3677 (2014).

835 75 Walker, L. C. *et al.* Evaluation of copy-number variants as modifiers of breast and  
836 ovarian cancer risk for BRCA1 pathogenic variant carriers. *Eur. J. Hum. Genet.* **25**,  
837 432-438 (2017).

838 76 Nelson, P. T., Fardo, D. W. & Katsumata, Y. The MUC6/AP2A2 Locus and Its  
839 Relevance to Alzheimer's Disease: A Review. *J. Neuropathol. Exp. Neurol.* **79**, 568-  
840 584 (2020).

841 77 Giannuzzi, G. *et al.* The Human-Specific BOLA2 Duplication Modifies Iron  
842 Homeostasis and Anemia Predisposition in Chromosome 16p11.2 Autism Individuals.  
843 *Am. J. Hum. Genet.* **105**, 947-958 (2019).

844 78 Mafrá, F. *et al.* Copy number variation analysis reveals additional variants  
845 contributing to endometriosis development. *J. Assist. Reprod. Genet.* **34**, 117-124  
846 (2017).

847 79 Jin, X. *et al.* Copy Number Variation of Immune-Related Genes and Their  
848 Association with Iodine in Adults with Autoimmune Thyroid Diseases. *Int. J.*  
849 *Endocrinol.* **2018**, 1705478 (2018).

850 80 Grau, C. *et al.* Xp11.22 deletions encompassing CENPVL1, CENPVL2, MAGED1  
851 and GSPT2 as a cause of syndromic X-linked intellectual disability. *PLoS One* **12**,  
852 e0175962 (2017).

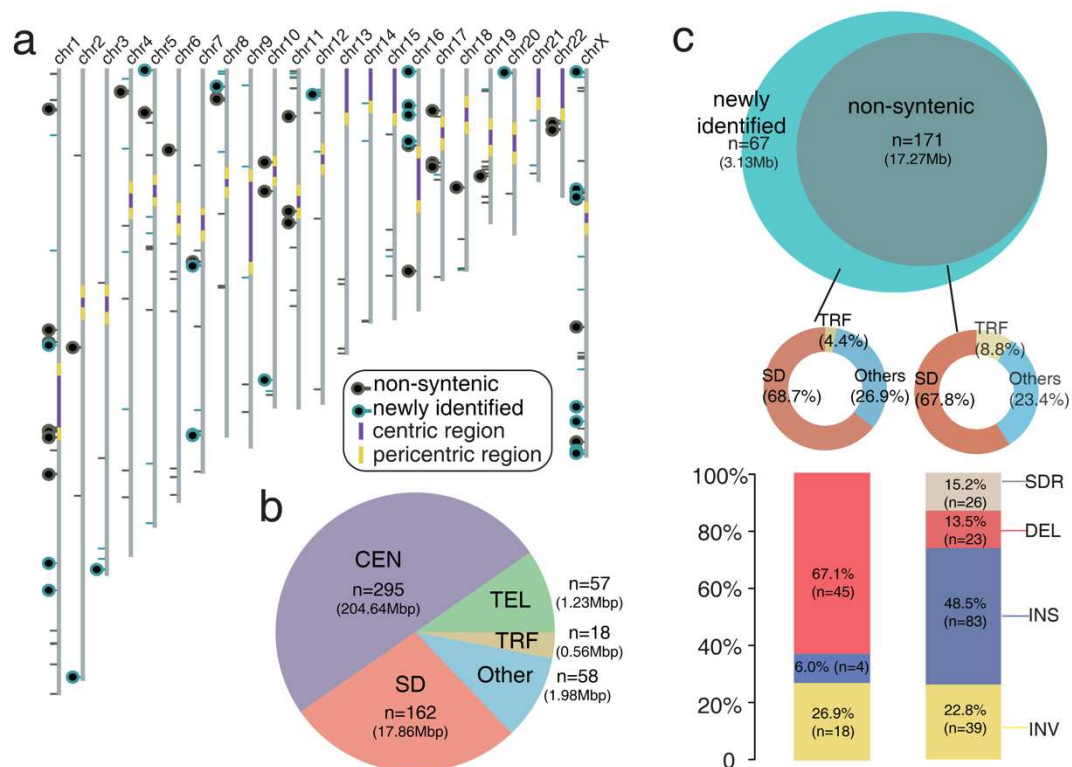
853 81 He, Y. *et al.* P2RY8 variants in lupus patients uncover a role for the receptor in  
854 immunological tolerance. *J. Exp. Med.* **219** (2022).

855 82 Wen, M. *et al.* CT45A1 promotes the metastasis of osteosarcoma cells in vitro and in  
856 vivo through  $\beta$ -catenin. *Cell Death Dis.* **12**, 650 (2021).

857 83 Yang, S. W. *et al.* A Cancer-Specific Ubiquitin Ligase Drives mRNA Alternative  
858 Polyadenylation by Ubiquitinating the mRNA 3' End Processing Complex. *Mol. Cell*  
859 **77**, 1206-1221.e1207 (2020).



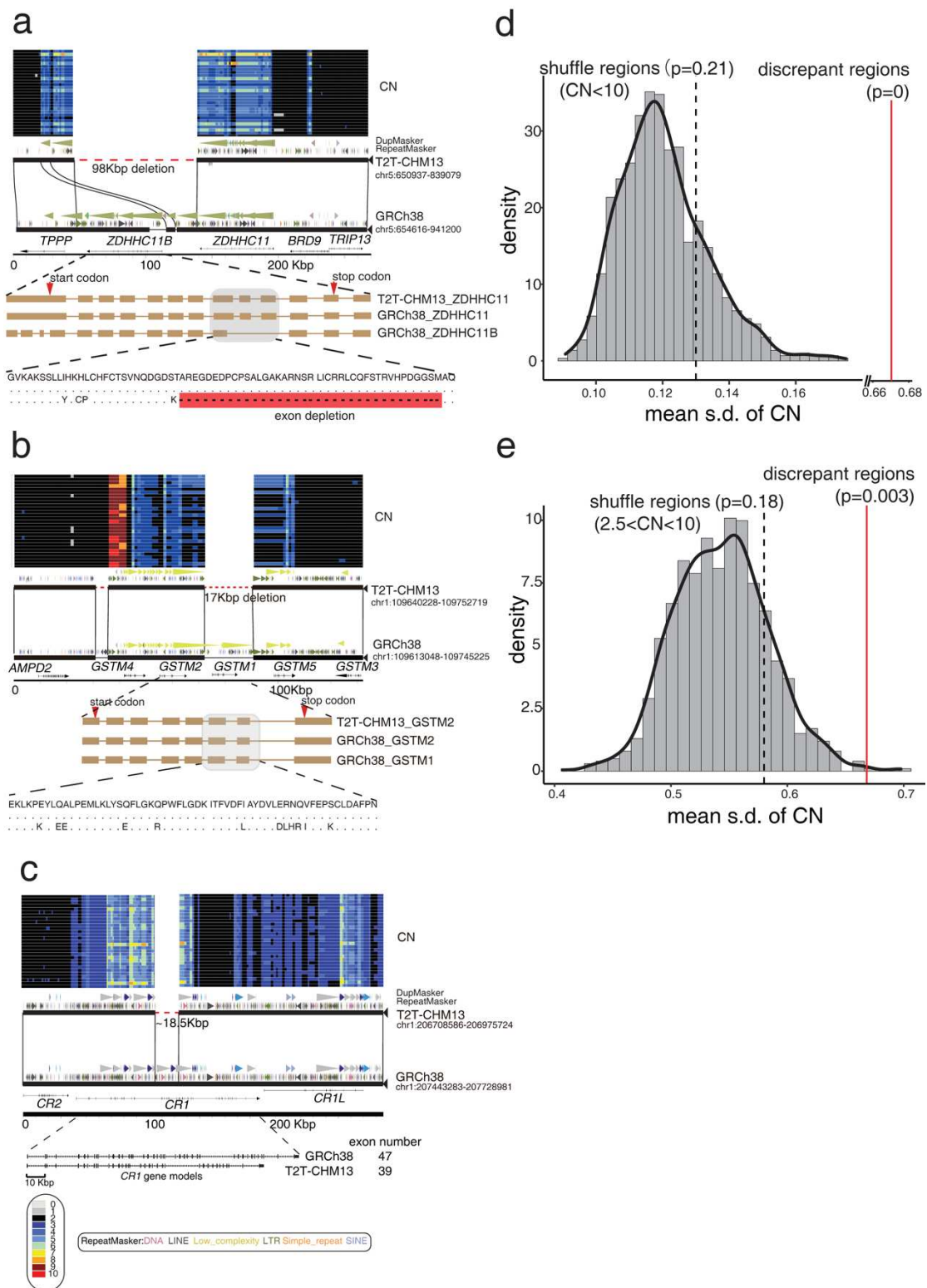
## Figure legends



**Figure 1. The discrepant genomic regions between GRCh38 and T2T-CHM13. (a)**

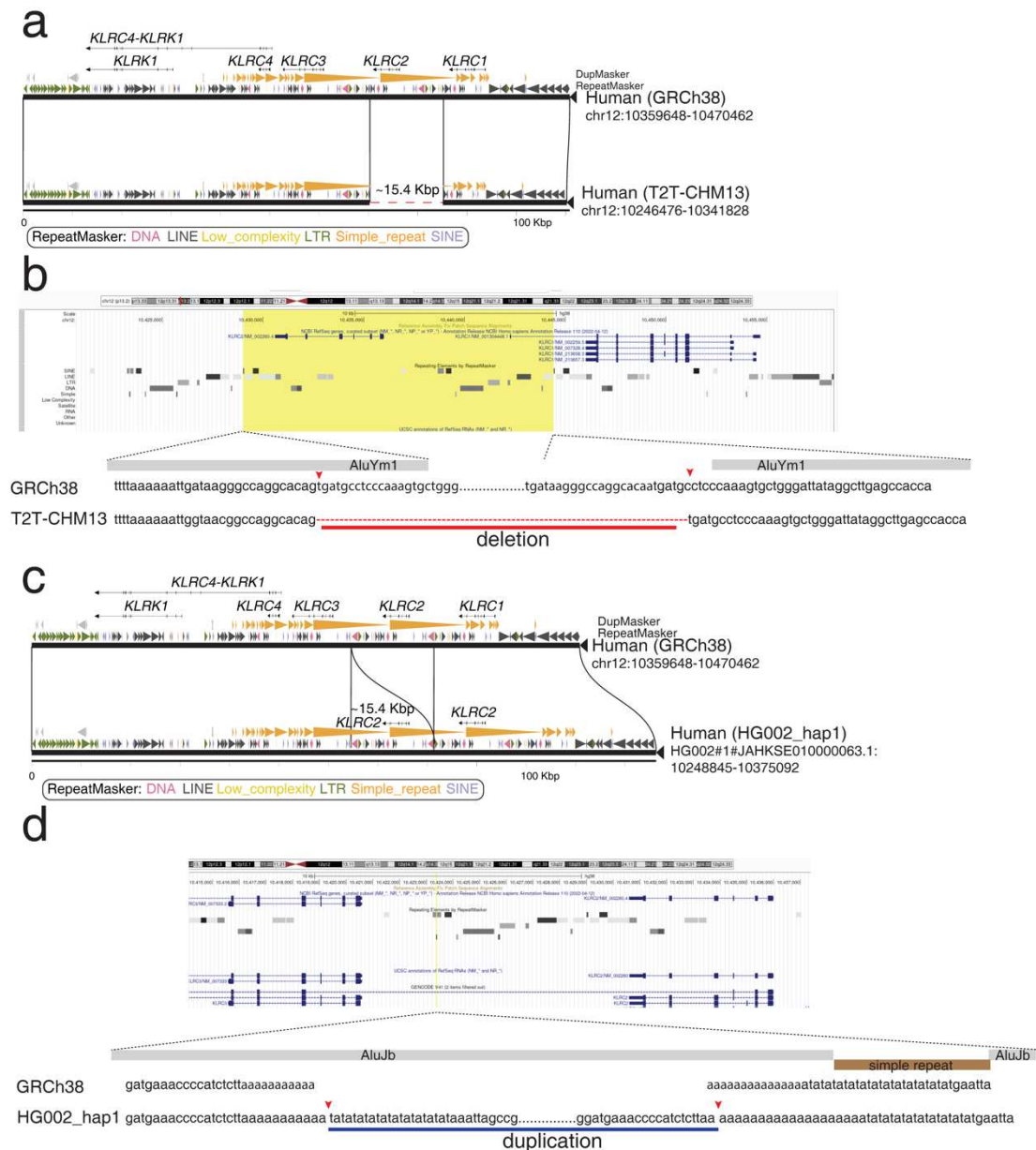
Schematic of the T2T-CHM13 assembly depicts the centromere location (purple and yellow), 'non-syntenic' region (black lines and circles), and newly identified discrepant region (cyan lines and circles). Regions containing genes are represented with circles. (b) Pie chart of genomic structure annotations of the 590 discrepant regions. The proportion of regions in centromeres (CEN), telomeres (TEL), segdups (SD), tandem repeats (TRF), and others are shown in light purple, green, dark red, yellow, and blue. (c) Venn diagram shows the comparison of the discrepant regions between the previous studies<sup>4,9</sup> and this study. The genome structure annotations of 'non-syntenic' and newly identified regions are shown in the middle panel. The components of structural variant types of 'non-syntenic' and newly identified regions are shown in the bottom panel.





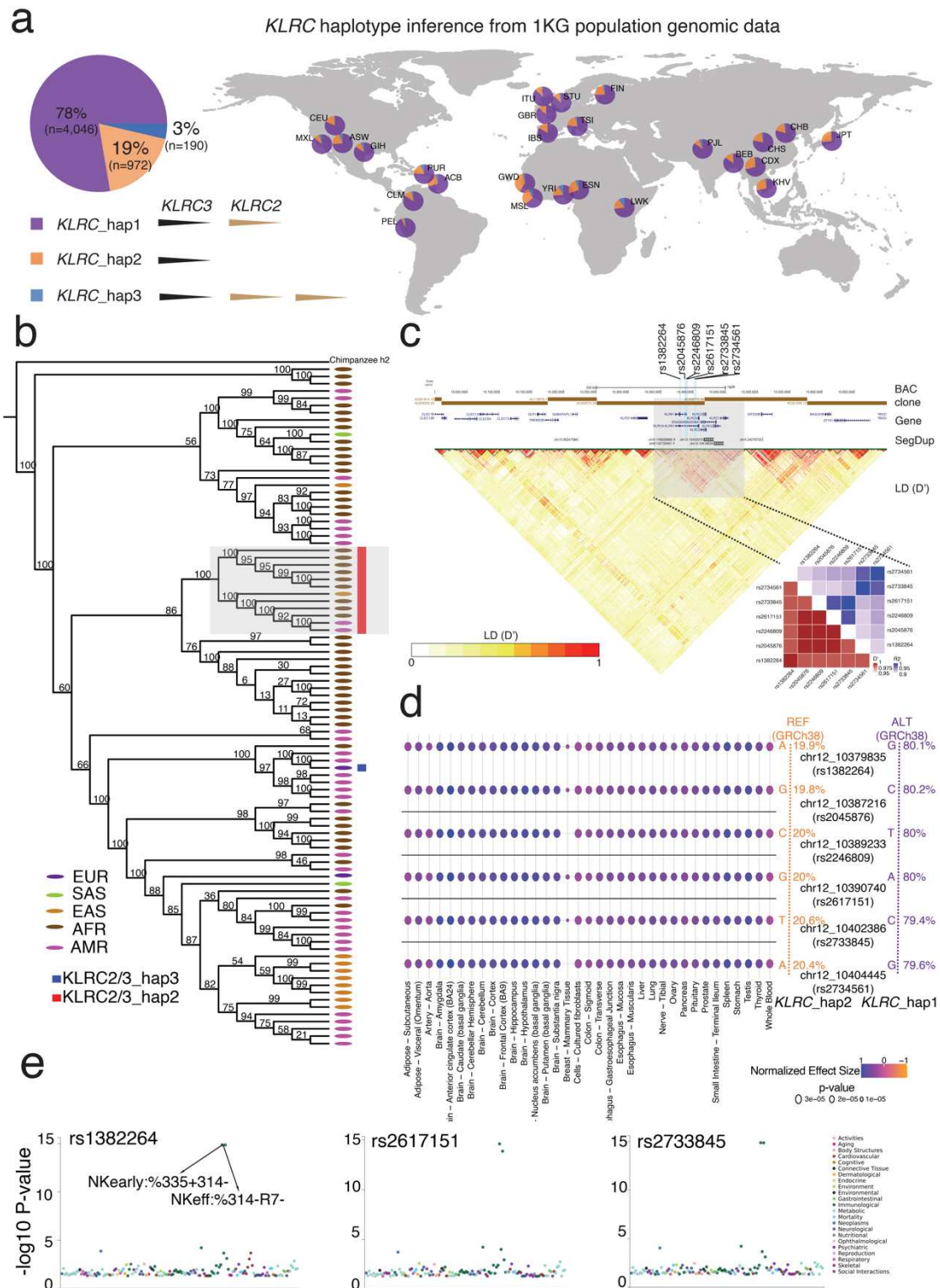
**Figure 2. Gene structure differences in the discrepant regions.** (a) The depletion of *ZDHHC11B* in the T2T-CHM13 genome assembly by a ~98 kbp deletion. The CN heatmap inferred from SGPD is shown in the top panel. The miropeat synteny relationship shows structural variation with repeat, duplication, and gene annotation. The exon schematic with amino acid alignment shows the gene model difference in the two assemblies. (b) The

880 depletion of *GSTM1* in the T2T-CHM13 genome assembly by a ~17 kbp deletion. (c) The  
 881 depletion of eight exons of *CRI* in T2T-CHM13 by ~18.5 kbp deletion. (d) The distribution  
 882 of the mean of s.d. of CN shows the mean s.d. of 131 discrepant regions (mean=0.735, red  
 883 line) is significantly higher than the simulated null distribution of s.d. of CN (CN<10,  
 884 empirical p=0). The black line represents the observed mean of s.d. of CN of the regions  
 885 where the CN is less than 10. (e) The distribution of the mean s.d. of CN shows the mean s.d.  
 886 of 131 discrepant regions (mean=0.735, red line) is significantly higher than the simulated  
 887 null distribution of s.d. of CN (2.5<CN<10, empirical p=0). The black line represents the  
 888 observed mean s.d. of CN of the regions where the CN is less than 10 and greater than 2.5.  
 889



**Figure 3. The syntenic comparison between different *KLRC* gene cluster haplotypes.** (a) A ~15.4 kbp deletion in the T2T-CHM13 genome assembly results in the complete loss of *KLRC2*. Gene structure, duplication, and repeat annotations are shown in the miropeat diagram. (b) A screenshot of the *KLRC2* region from the UCSC Genome Browser is shown in the top panel. The yellow box represents the ~15.4 kbp deleted sequence in the T2T-CHM13 genome assembly. The nuclear sequencing alignment of the breakpoints is shown in the bottom panel. Two Alu elements surrounding the breakpoints are shown in grey bars. (c) A ~15.4 kbp duplication in the HG002\_hap1 genome assembly results in the complete duplication of *KLRC2*. Gene structure, duplication, and repeat annotations are shown in the miropeats diagram. (d) A screenshot of the *KLRC2* region from the UCSC Genome Browser is shown in the top panel. The yellow line represents the position where the 15.4 kbp

902 duplicated sequence is in the HG002\_hap1 genome assembly. The nuclear sequencing  
903 alignment of the breakpoints is shown in the bottom panel. The disrupted Alu element within  
904 the breakpoints are shown in grey bars and a simple repeat disrupting the Alu element is  
905 shown in a brown bar.  
906

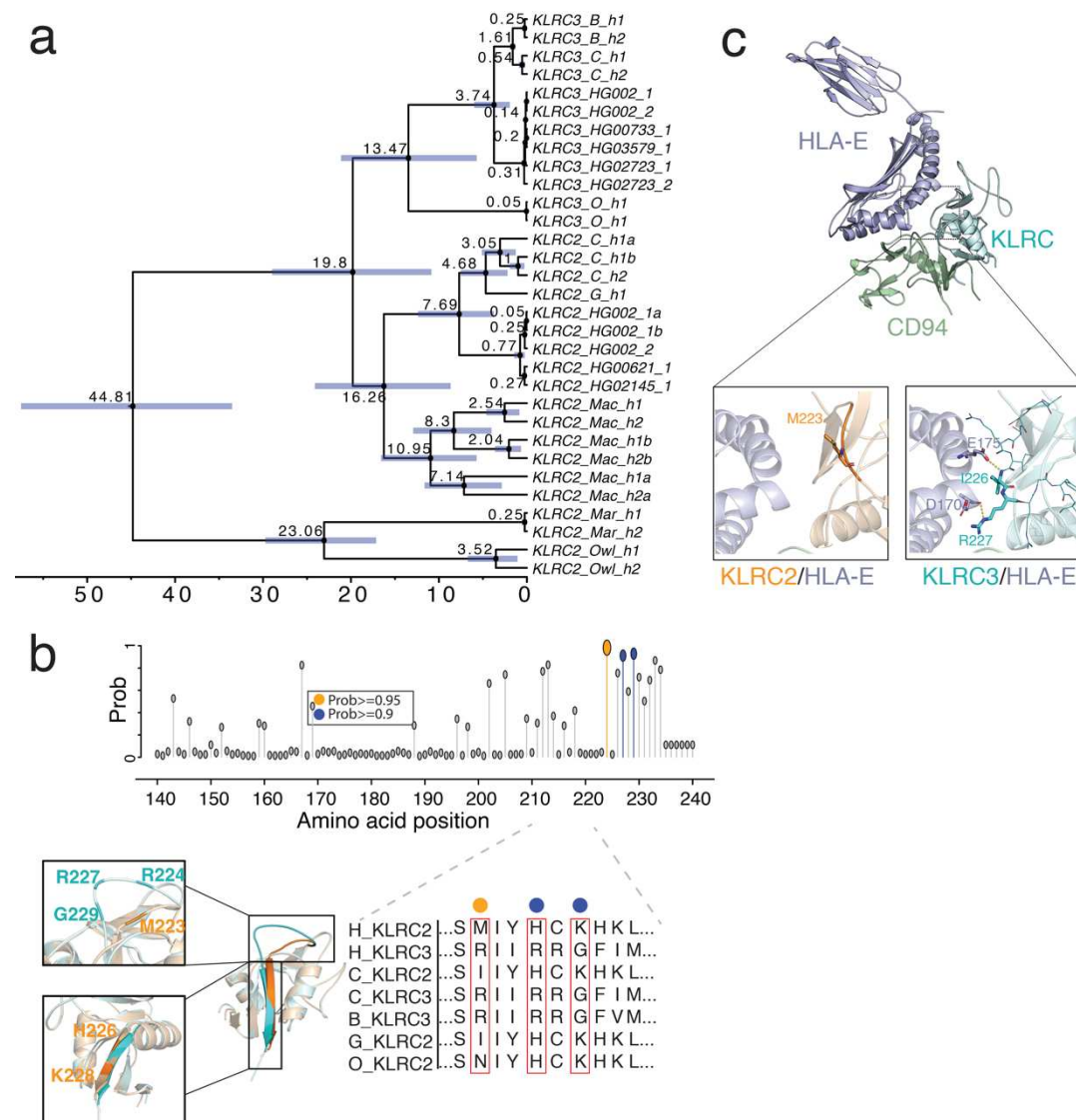


**Figure 4. The structural and functional diversity of *KLRC2* in humans.** (a) The proportion of three *KLRC* haplotypes is shown in a pie chart. The *KLRC\_hap1* represents one copy of *KLRC2* shown in purple. The *KLRC\_hap2* represents zero copies of *KLRC2* shown in orange. The *KLRC\_hap3* represents two copies of *KLRC2* shown in dark blue. Distributions of *KLRC\_hap1* (purple), *KLRC\_hap2* (orange) and *KLRC\_hap3* (blue) inferred from the 1KG



human population data across the world are shown on the right panel. (b) The phylogenetic tree of the *KLRC* haplotype genomic regions shows *KLRC*-hap2 (*KLRC2* depletion) is a result of a single deletion event. The red and blue rectangles show *KLRC*-hap2 and *KLRC*-hap3, respectively. The rest of the humans belong to *KLRC*-hap1. The super population of each human is listed with five color dots. (c) The genomic region (chr12:10,000,000-10,700,000 in GRCh38) with assembled BAC clone, gene, segdup, and LD (D') annotation shows that the *KLRC* gene cluster is probably linked together. The six SNVs distinguished between *KLRC*-hap2 and *KLRC*-hap1 are represented in the cyan lines with their SNP ID. The heatmaps of LD indexes (R2 and D') show that the six SNVs are highly linked in humans. (d) Consistent patterns of associations between the six SNVs and expression levels of *KLRC2* in 35 tissues are shown in the multi-tissue eQTL plots. The positive normalized effect size (NES) values represent the effect of the higher expression on the alternative allele (purple) relative to the reference allele (red). The (unadjusted) p values of the eQTL association are shown in the size variable dots. The allele frequencies of the six SNVs in the gnomAD database are shown on the right. (e) The PheWAS plots for three SNVs (rs1382264, rs2617151, and rs2733845) are significantly associated with immune domain differentiation across GWAS in the GWASALAS database. The particular traits (NKearly: %335+314- and NKeff: %314-R7-) are marked with significant signals<sup>32</sup>.





**Figure 5. The potential functional differentiation between *KLRC2* and *KLRC3* by natural selection in primate evolution.** (a) The phylogenetic tree reconstructed from *KLRC2* and *KLRC3* of humans and other NHPs with BEAST2 shows the duplication of *KLRC2* and *KLRC3* occurred at the common ancestor of African great apes. The 95% confidence interval of the estimated age of each node is shown in the blue bar. All nodes are supported by one posterior possibility shown in dark circle dots. The texts: C, B, G, O, Mac, Mar, and Owl in the tips represent chimpanzee, bonobo, gorilla, orangutan, macaque, marmoset, and owl monkey, respectively. (b) The possibility of amino acid under positive selection inferred by the branch-site model in PAML is shown on the top panel. The grey, orange, and blue dots represent the possibility of less than 0.9, between 0.9 and 0.95, or greater than 0.95, respectively. The amino acid alignment of *KLRC2* and *KLRC3* among primates is shown on the bottom right panel. Structure alignment of predicted structures of *KLRC2* from residue 118 to 231 (orange, Uniprot: P26717) and *KLRC3* from residue 118 to

946 240 (cyan, Uniprot: Q07444). The zoomed-in pictures depict the structural discrepancies in  
 947 the loop (top) and the following  $\beta$ -sheet (bottom) between KLRC2 and KLRC3. (c) Predicted  
 948 structures of KLRC/HLA-E/CD94 complex (KLRC2 from residue 118 to 231: orange,  
 949 Uniprot: P26717; KLRC3 from residue 118 to 240: cyan, Uniprot: Q07444; Full-length  
 950 HLA-E: purple, Uniprot: I3RW89; CD94 from residue 57 to 179: green, Uniprot: Q13241).  
 951 The zoomed-in protein structure depicts the interaction interfaces of KLRC2/HLA-E (top)  
 952 and KLRC3/HLA-E (bottom).

953 **Table 1. The discrepant regions associated with human diseases**

Cytobands	CHM13_Position	Hg38_Position	Type	Reported CNV	Genes	Disease
<b>1p13.3</b>	chr1:109711485-109711489	chr1:109682999-109701443	DEL		<i>GSTM1</i>	Urinary system disease <sup>36</sup>
<b>1p21.1</b>	chr1:103546781-103735057	chr1:103697900-103697950	INV		<i>AMY1A, AMY1B, AMYP1</i>	Neurological disease <sup>71</sup>
<b>1p36.13</b>	chr1:16007445-16027869	chr1:16565700-16565800	INS		<i>NBPF1</i>	Neurodevelopmental disorders <sup>72</sup> , Cancer
<b>1q21.1-1q21.2</b>	chr1:143959965-143983984	chr1:146251047-148716074	INV		<i>BCL9, NOTCH2NLB, CHD1L, NBPF12, PRKAB2, FMO5, ACP6, GJA8, GPR89B, NBPF11, NBPF14, PPIAL4G</i>	Neurodevelopmental disorders <sup>13-16</sup>
<b>1q31.3</b>	chr1:196105143-196105148	chr1:196758727-196843410	DEL		<i>CFHR3, CFHR1</i>	Immunological disease <sup>38</sup>
<b>1q32.2</b>	chr1:206810072-206810076	chr1:207542838-207561393	DEL		<i>CR1</i>	Neurological disease <sup>51,52</sup>
<b>2q13</b>	chr2:110517534-110698558	chr2:110095177-110276210	INV	Morbid CNV & Disease-related CNV	<i>NPHP1, MALL, MTLN</i>	Neurodevelopmental disorders <sup>28,29</sup> , Neurological disease <sup>30</sup>
<b>3q29</b>	chr3:198347865-198715835	chr3:195641035-195995576	DEL	Disease- related CNV	<i>MUC20, MUC4, TNK2</i>	Neurological disease <sup>30</sup>
<b>5p15.33</b>	chr5:684792-685093	chr5:686991-779053	DEL		<i>ZDHHC11B</i>	Cancer <sup>37</sup>
<b>6p21.32</b>	chr6:32339743-32356931	chr6:32486765-32530206	DEL		<i>HLA-DRB5</i>	Immunological disease <sup>73</sup>
<b>6q26</b>	chr6:161865491-161959834	chr6:160612509-160612509	INS		<i>LPA</i>	Cardiovascular disease <sup>74</sup>
<b>7q35</b>	chr7:145477647-145477649	chr7:144197172-144295737	DEL		<i>OR2A42, OR2A7, CTAGE8</i>	Cancer <sup>75</sup>
<b>8p23.1</b>	chr8:750030-11722000	chr8:8022351-12234558	INV	Morbid CNV	<i>DLGAP2, MYOM2, CLN8, ARHGEF, CSMD1, MCPH1, ANGPT2, PRR23D1, DEFB103B, DEFB103A, DEF104A, DEF105A, XKR6, SOX7, TNKS, PPP1R3B, PPAG1, CTSB, ANGPT2, AGPAT5, ERI1, MSRA, DEFA5, FDFIT1, GATA4, MFHAS1, PRSS5</i>	Developmental disorders <sup>28,29</sup>
<b>10q11.22</b>	chr10:48671598-48719249	chr10:47780140-47870155	SDR	Disease-related CNV	<i>GPRIN2</i>	Neurological disease <sup>31</sup>
<b>11p15.5</b>	chr11:1076897-1087865	chr11:1017980-1017990	INS		<i>MUC6</i>	Neurological disease <sup>76</sup>
<b>12p13.2</b>	chr12:10315804-10315827	chr12:10429009-10444430	DEL		<i>KLRC2</i>	COVID-19 <sup>34</sup> , Immunological disease <sup>32</sup>
<b>16p11.2</b>	chr16:30492288-30594258	chr16:30207700-30207750	INS	Disease- related CNV	<i>NPIPB13, BOLA2B</i>	Neurological disease <sup>30,77</sup>
<b>16p12.1-12.2</b>	chr16:28619710-29091966	chr16:28339205-28811381	INV	Disease- related CNV	<i>SULT1A1, SULT1A2, NPIPB8, NPIPB6, EIF3CL, NPIPB7, CLN3, IL27, EIF3C, NPIPB9</i>	Neurological disease <sup>30</sup>
<b>17p11.2</b>	chr17:16716173-16767175	chr17:16813513-16821452	SDR	Disease-related CNV	<i>LGALS9C</i>	Neurological disease <sup>30</sup>
<b>17q12</b>	chr17:37341285-37441106	chr17:36393230-36459266	SDR	Disease-related CNV	<i>CCL3L1, CCL4L2, TBC1D3F</i>	Neurological disease <sup>30</sup>
<b>19q13.2</b>	chr19:42710594-42726422	chr19:39906200-39906250	INS		<i>FCGBP</i>	Reproductive system disease <sup>78</sup>

<b>20p13</b>	chr20:1629529-1629530	chr20:1580346-1613395	DEL		<i>SIRPB1</i>	Immunological disease <sup>79</sup>
<b>22q11.23</b>	chr22:24380000-24462473	chr22:23932712-24000827	SDR	Disease-related CNV	<i>GSTT2, GSTT4, DDT</i>	Neurological disease <sup>30</sup>
<b>Xp11.22</b>	chrX:50939534-50996879	chrX:51668108-51725222	INV		<i>CENPVL1, CENPVL2</i>	Neurodevelopmental disorders <sup>80</sup>
<b>Xp22.33</b>	chrX:1307333-1307498	chrX:1465426-1506104	DEL		<i>P2RY8</i>	Immunological disease <sup>81</sup>
<b>Xq26.3</b>	chrX:134047172-134104452	chrX:13572163-135795043	SDR		<i>CT45A1, CT45A3, CT45A5</i>	Cancer <sup>82</sup>
<b>Xq28</b>	chrX:147946883-147987904	chrX:14968112-149722143	SDR		<i>MAGEA11</i>	Cancer <sup>83</sup>

954 Morbid CNV refers to the Ref. 28&29. Disease-related CNV refers to the Ref. 30.